

## Basic Statistics

732A47 Text Mining

## Overview

### Probability theory

- Probability
- Distributions and their properties
- Bayes theorem
- Generating test data

### Statistical Inference

- Hypothesis testing
- Linear regression
- Logistic regression
- Support vector machines

732A47 Text Mining

2

## Probability

How likely it is that some event will happen?

### Idea:

- Experiment
- Outcomes (sample points)  $O_1, O_2, \dots, O_n$
- Sample space  $\Omega$
- Event  $A$
- Probability function  $P$ : Events  $\rightarrow [0,1]$

**Example:** Tossing the coin

732A47 Text Mining

3

## Properties and definitions

- One can think of events as sets
  - Set operations are defined:  $A \cup B, A \cap B, \bar{A} \setminus B$
- $P(A \cup B) = P(A) + P(B)$  if  $A \cap B = \emptyset$
- **Independence**  $P(A, B) \equiv P(A \cap B) = P(A)P(B)$
- **Conditional probability**  $P(A|B) = \frac{P(A, B)}{P(B)}$

732A47 Text Mining

4

## Bayes theorem

### Example:

- We have constructed spam filter that
  - identifies spam mail as spam with probability 0.95
  - Identifies usual mail as spam with probability 0.005
- This kind of spam occurs once in 100,000 mails
- If we found that a letter is a spam, what is the probability that it is actually a spam?

732A47 Text Mining

5

## Bayes theorem

- We have some knowledge about event B
  - Prior probability  $P(B)$  of B
- We get new information A
  - $P(A)$
  - $P(A|B)$  probability of A can occur given B has occurred
- New (updated) knowledge about B
  - Posterior probability  $P(B|A)$

$$P(B|A) = \frac{P(A|B)P(B)}{P(A)}$$

732A47 Text Mining

6

## Bayes theorem for events

- If  $B_i$  are disjoint and  $A \subseteq \bigcup B_i$  then

$$P(B_j|A) = \frac{P(A|B_j)P(B_j)}{\sum_i P(A|B_i)P(B_i)}$$

Note,  $\bigcup B_i$  can be equal to  $\Omega$

732A47 Text Mining

7

## Random variables

- Instead of having events, we can have a variable X:
  - Events  $\rightarrow \mathbb{R}$  Continuous random variables
  - Events  $\rightarrow \mathbb{N}$  Discrete random variables

### Examples:

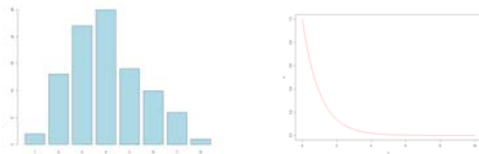
- $X = \{\text{amount of times the word "crisis" can be found in financial documents}\}$ 
  - $P(X=3)$
- $X = \{\text{Time to download a specific file to a specific computer}\}$ 
  - $P(X=0.36 \text{ min})$

732A47 Text Mining

8

## Distributions

- Discrete
  - Probability mass function  $P(x)$  for all feasible  $x$
- Continuous
  - Probability density function  $f(x)$
  - Cumulative density function  $F(x) = \int_0^x f(t)dt$



732A47 Text Mining

9

## Expected value and variance

- Expected value = mean value
  - $E(X) = \sum_{i=1}^n X_i P(X_i)$
  - $E(X) = \int X f(X) dX$
- Variance how much values of random variable can deviate from mean value
  - $Var(X) = E(X - E(X))^2 = E(X^2) - E(X)^2$

732A47 Text Mining

10

## Some conventional distributions

### Bernoulli distribution

- Events: Success ( $X=1$ ) and Failure ( $X=0$ )
- $P(X=1)=p$ ,  $P(X=0)=1-p$
- $E(X) = p$
- $Var(X) = 1 - p$

**Examples:** Tossing coin, winning a lottery,...

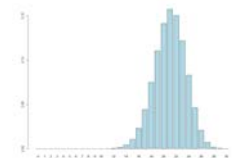
732A47 Text Mining

11

## Some conventional distributions

### Binomial distribution

- Sequence of  $n$  Bernoulli events
- $X = \{\text{Amount of successes among these events}\}$ ,  $X=0, \dots, n$
- $P(X=r) = \frac{n!}{(n-r)!r!} p^r (1-p)^{n-r}$
- $EX = np$
- $Var(X) = np(1-p)$



732A47 Text Mining

12

## Poisson distribution

- Customers of a bank  $n$  (in theory, endless population)
- Probability that a specific person will make a call to the bank between 13.00 and 14.00 a certain day is  $p$ 
  - $p$  can be very small if population is large (rare event)
  - Still, some people will make calls between 13.00 and 14.00 that day, and their amount may be quite big
  - A known quantity  $\lambda=np$  is mean amount of persons that call between 13.00 and 14.00
  - $X=\{\text{amount of persons that have called between 13.00 and 14.00}\}$

732A47 Text Mining

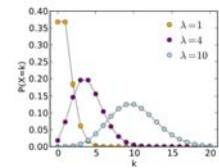
13

## Poisson distribution

- $P(X = r) = \lim_{n \rightarrow \infty} \frac{n!}{(n-r)!r!} p^r (1-p)^{n-r}$
- It can be shown that

$$P(X = r) = \frac{\lambda^r e^{-\lambda}}{r!}$$

- $E(X) = \lambda$
- $Var(X) = \lambda$



732A47 Text Mining

14

## Poisson distribution

- Further properties:
  - Poisson distribution is a good approximation of the binomial distribution if  $n > 20$  and  $p < 0.05$
  - Excellent approximation if  $n \geq 100$  and  $np \leq 10$

732A47 Text Mining

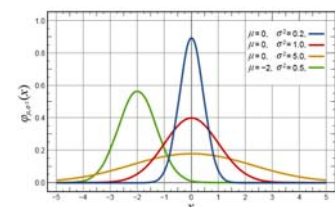
15

## Normal distribution

- Appears in almost all applications
  - Time required to download a specific document to a specific computer

$$f(x) = \frac{1}{\sqrt{2\pi}\sigma} e^{-\frac{(x-\mu)^2}{2\sigma^2}}, \sigma > 0$$

- $E(X) = \mu$
- $Var(X) = \sigma^2$

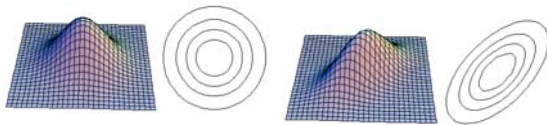


732A47 Text Mining

16

## Multivariate distributions

- Probability of two variables having certain values at the same time
  - P.D.F.  $f(x,y)$
  - Correlation



732A47 Text Mining

17

## Radnom number generation

- Random number (sequences) can be generated in computer but those are not truly random!
  - Deterministic algorithms are used
  - Current time is also used to mimic randomness

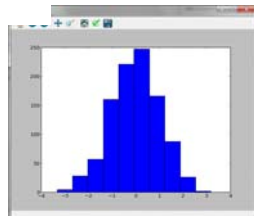
```
>>> from scipy.stats import norm
>>> norm.pdf([-1., 0, 1])
array([ 0.24197072, 0.39894228, 0.24197072])
>>> norm.pdf(0, loc=1, scale=2)
0.17603266338214976
>>> t=norm.rvs(loc=2,scale=0.1,size=10)
>>> t
array([ 2.12425886, 1.81937297, 2.15002352, 2.086603 , 2.1119639 ,
        1.8902981 , 1.98702357, 2.03978701, 2.00005928, 2.13550148])
>>>
```

732A47 Text Mining

18

## Generating a sample

```
>>> from scipy.stats import norm
>>> t=norm.rvs(size=1000)
>>> hist(t)
```



732A47 Text Mining

19

## Generating datasets

- Necessary for testing different algorithms
- Dataset  $D = \{(X_i, Y_i), i = 1, \dots, n\}$

Case	Variable 1	Variable 2	Variable 3	.	.	.	.
1							
2							
3							
.							
.							
.							

- Decide
  - the structure of each input variable
  - the distribution of some input variables
  - multivariate distribution for some subset of input variables
  - Response model

732A47 Text Mining

20

## Generating data

- Response models
  - Continuous variables
    - $Y = f(X)$  or  $y = f(X) + \epsilon$
  - Discrete variables
    - $Y = f(X)$  or  $y = \text{Generator}_p(f(X))$

732A47 Text Mining

21

## Hypothesis testing

- Certain proportion of population satisfies some property in our data,  $p$
- The true distribution is *Bernoulli*( $\pi$ )
- We have assumption about  $\pi$ :
  - Test  $\pi = \pi_0$  vs  $\pi \neq \pi_0$
- Special statistical procedures exist to find it out
  - Compute  $s = \sqrt{p(1-p)}$
  - Compute the absolute value of  $z = \frac{p-\pi_0}{\sqrt{\frac{s^2}{n}}}$  and compare it with value from a table
  - If this value is greater the table value, reject  $\mu = 0.5$

732A47 Text Mining

22

## Chi-square test

H0: Frequencies of words are approximately same in all documents

Ha: Some documents have a different pattern

Word	Document1	Document 2
business	43	66
school	38	72
dollar	11	23
market	8	19

$$\chi^2 = \sum \frac{(O-E)^2}{E}$$

732A47 Text Mining

23

## Simple linear regression

**Model:**

$$y = \beta_0 + \beta_1 x + \text{error}$$

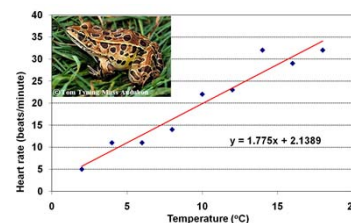
**Terminology:**

$\beta_0$ : intercept (or bias)

$\beta_1$ : regression coefficient

**Response**

The target responds directly and linearly to changes in the input



Data mining and statistical learning-2012

## Ordinary least squares regression (OLS)

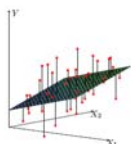
### Model:

$$y = \beta_0 + \beta_1 x_1 + \dots + \beta_p x_p + \text{error}$$

$$y = \beta_0 + \beta^T \mathbf{X} + \text{error}$$

Find coefficients by minimizing RSS:

$$RSS(\beta) = \sum_{i=1}^N (y_i - \beta_0 - \beta_j \sum_{j=1}^p x_{ij})^2$$



Data mining and statistical learning-2012

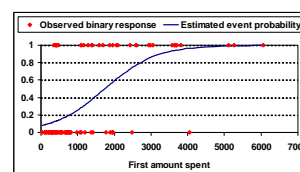
## Logistic regression: two classes

- Consider Logistic model with predictor  $X = \text{First\_Amount\_Spent}$

### Logistic model

$$\log \frac{P(Y=1|X=x)}{P(Y=0|X=x)} = \log \frac{P(Y=1|X=x)}{1-P(Y=1|X=x)} = \beta_0 + \beta_1 x$$

$$P(Y=1|X=x) = \frac{\exp(\beta_0 + \beta_1 x)}{1 + \exp(\beta_0 + \beta_1 x)}$$



732A33-2012

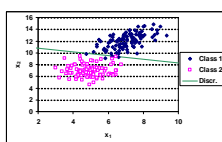
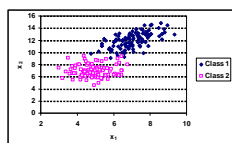
26

## Classification

Given a dataset  $D = \{(\mathbf{X}_i, Y_i), i = 1, \dots, n\}$

- $X_1$ =frequency of "money"
- $X_2$ =frequency of "crisis"
- $Y=0$  if document is from financial area
- $Y=1$  if document is from "psychology"

Aim: divide input space into areas in order to predict new observations



732A47 Text Mining

27

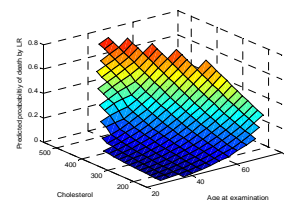
## Logistic regression: two classes

- Two predictors

— Fitted probability of death in coronary heart disease as a function of cholesterol level and age at examination

$$\log \frac{P(Y=1|X=\mathbf{x})}{P(Y=0|X=\mathbf{x})} = \beta_0 + \beta_1^T \mathbf{x}$$

$$P(Y=1|X=\mathbf{x}) = \frac{\exp(\beta_0 + \beta_1^T \mathbf{x})}{1 + \exp(\beta_0 + \beta_1^T \mathbf{x})}$$



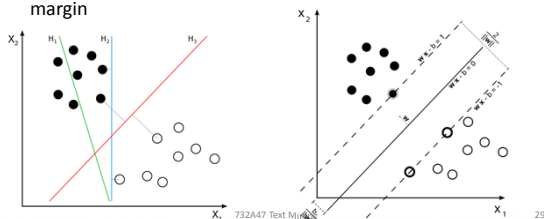
732A33-2012

28



## Support vector machines

- Very good in classification of high dimensional data → very popular for text classification
- Support vector=observations on the boundary of classes
- **Linear SVM model**: find a linear boundary having maximal margin



732A47 Text Mining 29

## Linear SVM

- Boundary is given by equation:

$$-w \cdot x + b = 0$$

- Condition for  $Y=1$ :

$$-w \cdot x + b \geq 1$$

- Condition for  $Y=-1$ :

$$-w \cdot x + b \leq -1$$

$$\min_w \|w\|^2 \quad \text{subject to } y_i(w \cdot x_i + b) \geq 1, \quad i = 1, \dots, n$$

→ Convex optimization problem

We need to find  $w$  and  $b$  that have largest margin  $\frac{2}{\|w\|}$

732A47 Text Mining

30

## SVM

- Compared to many other methods, a global minimum of the objective function is possible to find
- There are generalizations of SVM for
  - Nonlinear boundaries
  - More than two classes
  - Cases when the boundaries are not separable

732A47 Text Mining

31

## Reading

- [www.scipy.org](http://www.scipy.org)
- MS, Chapter 2, 5.1-5.3

732A47 Text Mining

32